

Cite this: *Mater. Horiz.*, 2025, 12, 7416Received 14th April 2025,  
Accepted 25th June 2025

DOI: 10.1039/d5mh00699f

rsc.li/materials-horizons

# Deep learning-enhanced development of innovative antioxidant liposomal drug delivery systems from natural herbs†

Xiaohe Zhang,<sup>‡a</sup> Zhihang Zheng,<sup>‡a</sup> Lina Xie,<sup>b</sup> Minghao Yang,<sup>a</sup> Jing Wang,<sup>c</sup> Weiwei Wang,<sup>c</sup> Shuyan Han,<sup>\*d</sup> Zhen Zhang<sup>\*b</sup> and Jun Wu<sup>‡\*ae</sup>

Free radical-mediated oxidative damage to biological macromolecules, such as DNA and proteins, significantly contributes to cellular ageing. Antioxidants play a crucial role in mitigating this process by neutralizing reactive oxygen species (ROS) and reducing DNA damage. Traditional herbal medicines are of strong interest as potential sources of antioxidants due to their rich diversity of bioactive components. In this study, we developed a two-stage BERT-based framework trained on 587 experimentally confirmed antioxidants and 983 inactive compounds. The optimized model effectively screened a broad range of potential antioxidant compounds from a library of 2882 natural herbal compounds, achieving an accuracy improvement of approximately 20% over traditional machine learning models. Molecular docking simulations and *in vitro* experiments consistently validated the antioxidant capacity of the selected compounds. Additionally, incorporating three representative compounds into a liposomal delivery system not only enhanced *in vivo* bioavailability, but also mitigated oxidative stress injury after kidney acute ischemia/reperfusion. This was achieved by up-regulating antioxidant-related genes in target organs as well as ROS scavenging. Our findings highlight the potential of integrating deep learning-based compound screening with an engineered liposomal delivery platform in the research of oxidative stress and aging.

## 1. Introduction

Oxidative stress is one of the core pathological mechanisms of many chronic diseases (such as cancer, neurodegenerative and

### New concepts

This study constructs two stage framework architecture based on a pre-trained BERT model. The *t*-SNE analysis revealed that the model can spontaneously capture the clustering characteristics of chemical functional groups through extensive unsupervised pre-training, indicating that it has basic chemical semantic understanding capabilities. After fine-tuning the dataset provided by this study, the results show that its five-fold cross-validation average AUC value reaches 0.9832 and the accuracy rate reaches 0.9363, which is significantly better than traditional models (random forest, SVM) and CNN models. Through the analysis of the attention mechanism, the association between key molecular substructures such as the adjacent dihydroxy structure and biological activity was successfully identified, providing a structure-guided mechanistic explanation for drug design. The compounds screened from 2882 natural herbs based on the model were efficiently delivered by liposome modification technology, and the antioxidant properties of the virtual drug screening were verified in multiple dimensions in *in vitro* and *in vivo* experiments. This work provides a new paradigm for large-scale drug screening and design by fusing virtual screening with efficient delivery technology, while expanding the application scenarios of virtual drug screening in the fields of oxidative stress regulation, materials science and biomedicine.

cardiovascular diseases) and the ageing process.<sup>1–3</sup> Its essence is biomolecular damage caused by an imbalance between production and scavenging of free radicals in the body. Natural herbs, an important carrier of traditional medicine, are rich in antioxidant active ingredients (such as polyphenols, flavonoids, terpenes and alkaloids).<sup>4</sup> They have the unique advantages of multi-target regulation of redox balance, low toxicity and side

<sup>a</sup> Bioscience and Biomedical Engineering Thrust, The Hong Kong University of Science and Technology (Guangzhou), Nansha, Guangzhou 511400, China.

E-mail: junwuhkust@ust.hk

<sup>b</sup> Department of Hematology, The Seventh Affiliated Hospital, Sun Yat-sen University, Shenzhen 518107, China. E-mail: zhangzhen1@sysush.com

<sup>c</sup> Department of General Surgery, The First Affiliated Hospital of Sun Yat-sen University, Guangzhou, 510080, China

<sup>d</sup> Department of Nephrology, Center of Kidney and Urology, The Seventh Affiliated Hospital, Sun Yat-Sen University, Shenzhen 518107, China.

E-mail: hanshy26@mail.sysu.edu.cn

<sup>e</sup> Division of Life Science, The Hong Kong University of Science and Technology, 999077, Hong Kong

† Electronic supplementary information (ESI) available. See DOI: <https://doi.org/10.1039/d5mh00699f>

‡ These authors contributed equally to the article.



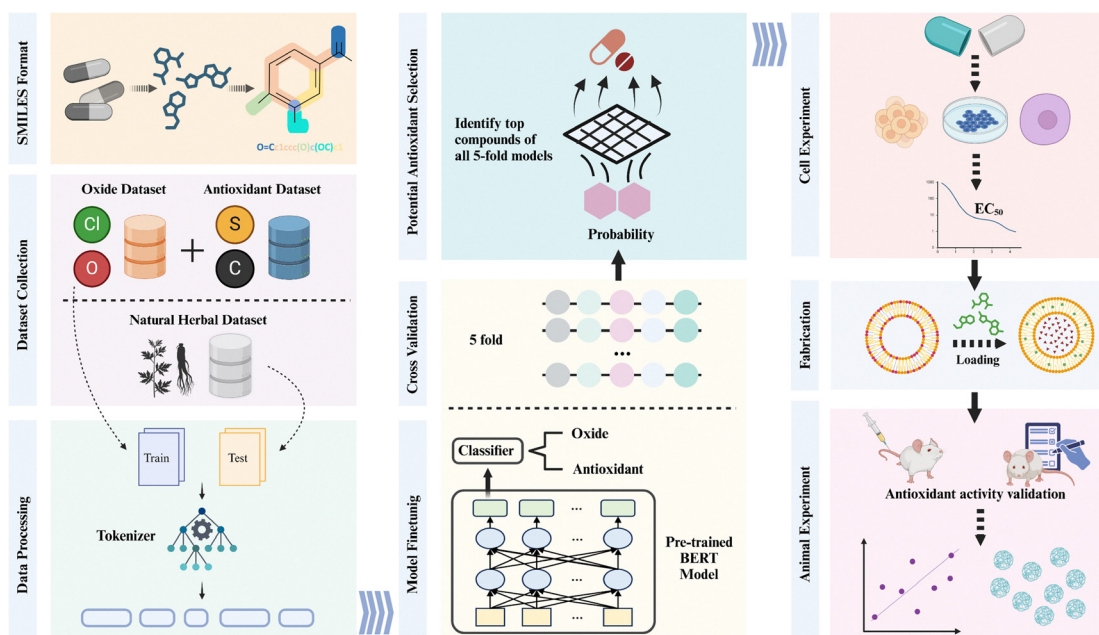
effects, and a wide range of sources. However, the main process for the discovery of active ingredients in traditional herbal is empirical screening and isolation, followed by purification and validation at cellular and animal levels. This process has problems such as low efficiency, high cost, high failure rate and difficulties in analyzing the synergistic effects of complex ingredients, which limits the process of modernizing their development.<sup>5</sup>

In recent years, the paradigm of natural product research and development has been transformed by breakthroughs in artificial intelligence (AI) technology, particularly deep learning algorithms. The drug discovery cycle has been significantly shortened by integrating multiple machine learning models and neural network architectures to discover new compounds.<sup>6–9</sup> Specifically, convolutional neural networks (CNNs), graph convolutional networks (GCNs), *etc.*, analyze the combination of high-order features between molecular connections to predict the properties of corresponding molecules. The transfer learning framework and active learning strategies can then efficiently analyze massive drug or compound databases to predict the biofunctionality of novel unknown compounds and their targets to power downstream drug development. The universality and convenience of such methods have been reported, for example, Stokes *et al.* used deep neural network (DNN) to predict a new antibiotic, halicin, and identified eight antibacterial compounds with large structural differences from known antibiotics from a database of more than 107 million molecules, greatly improving the efficiency of antibiotic library expansion.<sup>10</sup>

However, natural antioxidant ingredients generally have bottlenecks such as poor water solubility, low bioavailability and insufficient stability in the body, which limit their clinical application. As a novel drug delivery system, liposomes can simultaneously encapsulate hydrophilic and hydrophobic

active molecules due to their amphiphilic phospholipid bilayer structure and achieve long-lasting circulation, tissue-specific delivery and enhanced transmembrane penetration through surface modification (such as PEGylation and targeted ligands). In addition, liposomes can protect active ingredients from enzymatic or pH degradation and prolong the antioxidant effect through a sustained release mechanism.<sup>11–14</sup> For example, Liao *et al.* were inspired by traditional Chinese medicine compounds and developed bergamot liposomes for the treatment of acute respiratory distress syndrome (ARDS). Compared to free drug, the bioavailability was increased by almost 10-fold and there was significant targeting to the site of inflammation in the lung, significantly enhancing the efficacy of bergamot and reducing its systemic toxicity.<sup>15</sup>

In this study, we developed a BERT-based molecular antioxidant property prediction model and applied it to identify potential candidate compounds from natural herbs with about 20% higher accuracy compared to conventional machine learning models (RF and SVM). Specifically, this method uses a transformer architecture to successfully capture the underlying features of the SMILES structure of antioxidant compounds and autonomously learns to screen for novel compounds with potential antioxidant properties in a natural herbal compound library. Finally, the physicochemical properties of the newly discovered compounds are used to construct functional liposomes to verify their efficient antioxidant properties *in vitro* and *in vivo* (Scheme 1). Our research pioneering integrates deep learning and liposome technology to successfully construct an integrated platform of “intelligent prediction and efficient delivery”, providing a new paradigm for breaking through the barriers of natural ingredient delivery and developing safe and effective antioxidant treatment strategies. This cross-integration strategy not only promotes the modern use of natural herbal resources but also



Scheme 1 Illustration of BERT-based antioxidant molecular discovery model and antioxidant liposome construction and *in vitro* and *in vivo* validation.



opens new ideas for precision medicine in antioxidant protection and personalized disease treatment, while also opening innovative approaches for the development of other functional liposomes.

## 2. Results and discussion

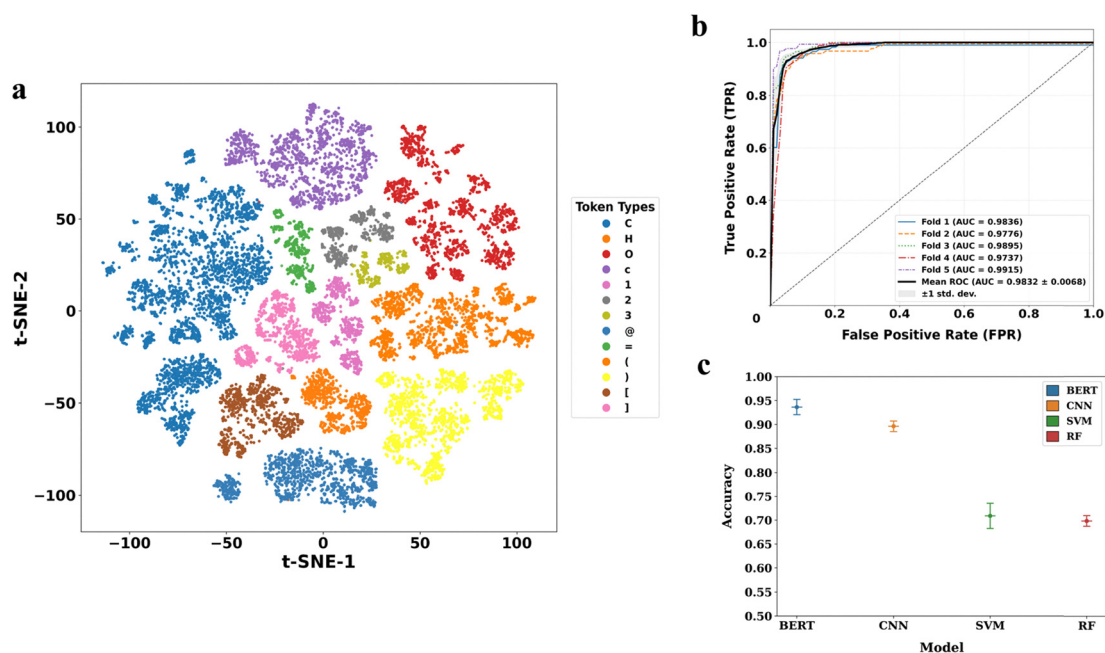
### 2.1 Training and optimization results for the antioxidant-related model

To elucidate the learned representations of the pre-trained BERT model for molecular structures, we implemented a pipeline focusing on token embedding distributions. 500 molecules were randomly sampled from the dataset, and embeddings were extracted for all constituent tokens. To focus on statistically significant patterns and improve visual clarity, frequency-based filtering was applied for the *t*-SNE (*t*-distributed stochastic neighbor embedding) visualization. Fig. 1a presents the *t*-SNE visualization for tokens with  $\geq 500$  occurrences, with comprehensive visualizations across multiple frequency thresholds provided in Fig. S1a–d (ESI<sup>†</sup>). As shown in Fig. 1a, there is a pronounced clustering phenomena where chemically equivalent functional groups consistently aggregated in the latent space despite originating from diverse molecular environments. This result demonstrates that after pre-training on large scale data, the model has acquired fundamental chemical understanding without explicit incorporation of molecular theory.

To identify compounds with antioxidant activity, we fine-tuned the pre-trained BERT model on our dataset containing 584 compounds with antioxidant activity and 983 with oxidizing

activity. The classification performance of model was evaluated using five-fold cross-validation. Fig. 1b presents the receiver operating characteristic (ROC) curves and area under the curve (AUC) values for each fold, along with the mean values. The mean AUC of 0.9832 demonstrates the fine-tuned model's strong ability to distinguish between antioxidants and oxidizers. We also compare the BERT model against three baseline approaches: random forest (RF), support vector machine (SVM), and convolutional neural network (CNN). The RF classifier leverages Morgan fingerprints with a radius of 2 and 2048 bits to capture circular substructures around each atom. The SVM model employs a comprehensive set of 25 physicochemical descriptors including molecular weight, hydrogen bond features, and topological properties. The CNN architecture encodes SMILES strings into 42-dimensional feature vectors—21 dimensions representing atomic properties and 21 dimensions employing one-hot encoding for chemical symbols—which are processed through multiple convolutional and pooling layers to generate SMILES Convolutional Fingerprints (SCFPs) for molecular classification.<sup>16</sup> Fig. 1c shows the five-fold cross-validation accuracy results across all models.

The BERT model outperforms others, which achieves the highest values for mean accuracy (0.9363). The CNN model, with a mean accuracy of 0.8962, is slightly worse than the BERT model. Compared with traditional machine learning models (RF and SVM), both deep learning models show obvious performance improvement with around 20% increase in accuracy. In short, the BERT model's superior performance stems from its ability to capture contextual relationships within SMILES representations. The result underscores the potential of transformer-based architectures in molecular property



**Fig. 1** Performance evaluation of the BERT model. (a) *t*-SNE visualization of representative molecular embeddings generated by the pre-trained BERT model, illustrating the correspondence between SMILES tokens, their embedding positions, and associated molecular atoms with consistent color coding. (b) ROC curves displaying the 5-fold cross-validation results, where the black line represents the mean performance with an average AUC of  $0.9832 \pm 0.0068$ . (c) Performance comparison between the BERT-based approach and benchmark models.



prediction, where the extensive parameterization and self-attention mechanisms effectively learn the complex structure–property relationships of antioxidant/oxidizer.

The BERT model used here leverages attention mechanisms to integrate information from SMILES tokens into task-specific molecular representations. Therefore, we examine the interpretation of molecular structure–property relationships using BERT attention mechanisms. This method addresses the need for interpretable molecular modeling by establishing direct correlations between structural elements and predicted properties. By analyzing attention patterns across SMILES tokens, salient molecular substructures contributing to predictions are identified, providing insights into structure–activity relationships that inform rational molecular design.

To validate the biological relevance of these attention assignments, we examine molecules from the natural herb set. Here, the attention scores from BERT tokens are mapped to corresponding characters in SMILES strings, with normalization applied to emphasize relative importance across molecular substructures. These character-level attention values are then projected onto 2D molecular representations using color gradient, where deeper indicates higher attention weights. Four bioactive molecules are selected for attention visualization analysis of their chemical structures and SMILES representations: two are free of antioxidant properties (Ascaridole and Artesunate) and two with confirmed antioxidant activities (Butein and Protocatechuic acid).<sup>17,18</sup> In Fig. 2a, it is shown that regions with higher attention (dark red) often represent aromatic ring conjugated system and *ortho*-dihydroxy structure. The aromatic rings provide stable electron conjugation systems, which stabilize the free radical intermediates formed after capturing free radicals, improving antioxidant efficiency.

While for the other two molecules in Fig. 2b, we can notice more attention being assigned to peroxide bridge, which are significant for oxidant ability. These findings demonstrate that BERT effectively assigns property-specific attention weights, offering medicinal chemists' valuable tools to explore connections between molecular substructures and their associated properties.

## 2.2 Molecular docking to verify antioxidant properties

To further verify whether the compounds screened by the model have antioxidant effect, we molecularly docked the key protein targets of antioxidant multiple pathways with them (Nrf2, SOD, and HO-1), to preliminarily evaluate the intensity of their endogenous antioxidant capacity. The above protein targets coordinate Nrf2/ARE, FOXO, SIRT1 *etc.* signaling pathways to construct the endogenous antioxidant network of the organism, which reduces the damage of DNA, proteins and lipids by oxidative stress, and delay related diseases process. Among them, related studies showed that fisetin and honokiol can activate the Nrf2 signaling pathway and promote the expression of antioxidant enzymes, which further backed up the accuracy of our machine model.<sup>19,20</sup> Molecular docking results showed that 5,7-diacetoxy-8-methoxyflavone, fisetin and honokiol can all stably bind to the active pocket of Nrf2. Visualization results show that 5,7-diacetoxy-8-methoxyflavone stably binds to Nrf2 through a hydrogen bond with a specific residue Arg-326 (distance: 2.4 Å) and a hydrophobic interaction with Gly-371 (distance: 1.8 Å). In addition, fisetin forms hydrogen bonds with the Nrf2 residues Asp-479, His-436 and three hydrogen bonds with Gly-480. Honokiol and Nrf2 residues Val-606 and Val-467 form two hydrogen bonds (Fig. 3a–c). Fig. 3d–f demonstrate that the flavonoid derivatives 5,7-diacetoxy-8-methoxyflavone, fisetin, and honokiol exhibit

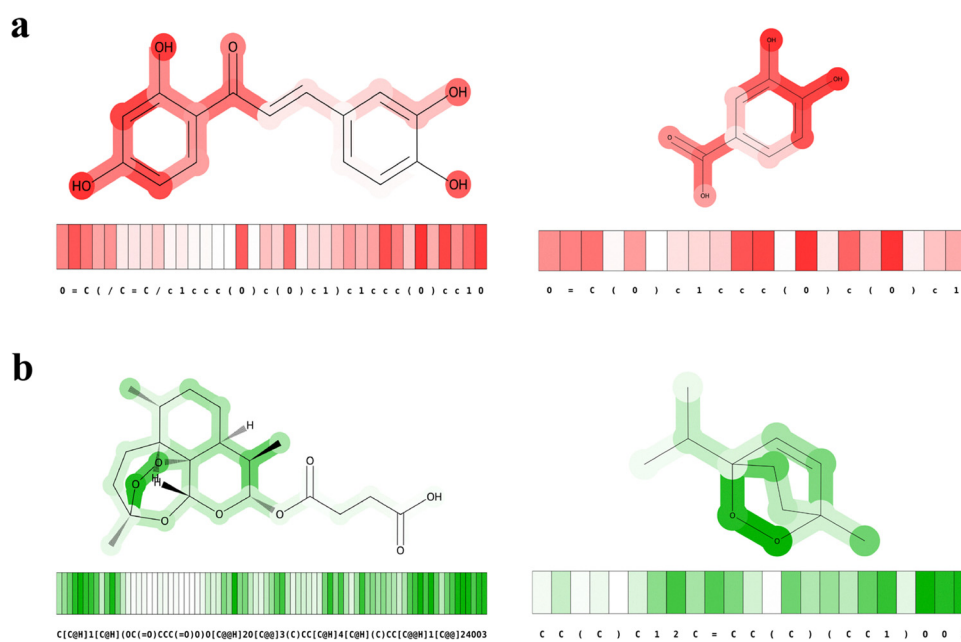
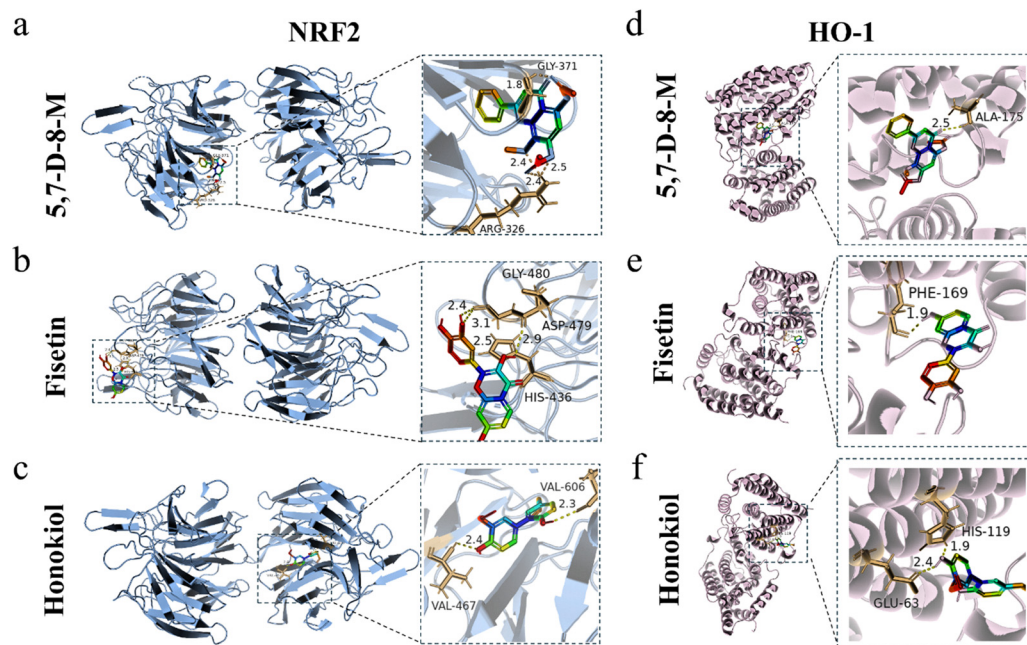


Fig. 2 Attention heat map of molecules with (a) antioxidant activity and (b) molecules with oxidant activity. BERT's attention scores are allocated to both SMILES tokens and structural representation. Darker colors indicate higher attention scores.

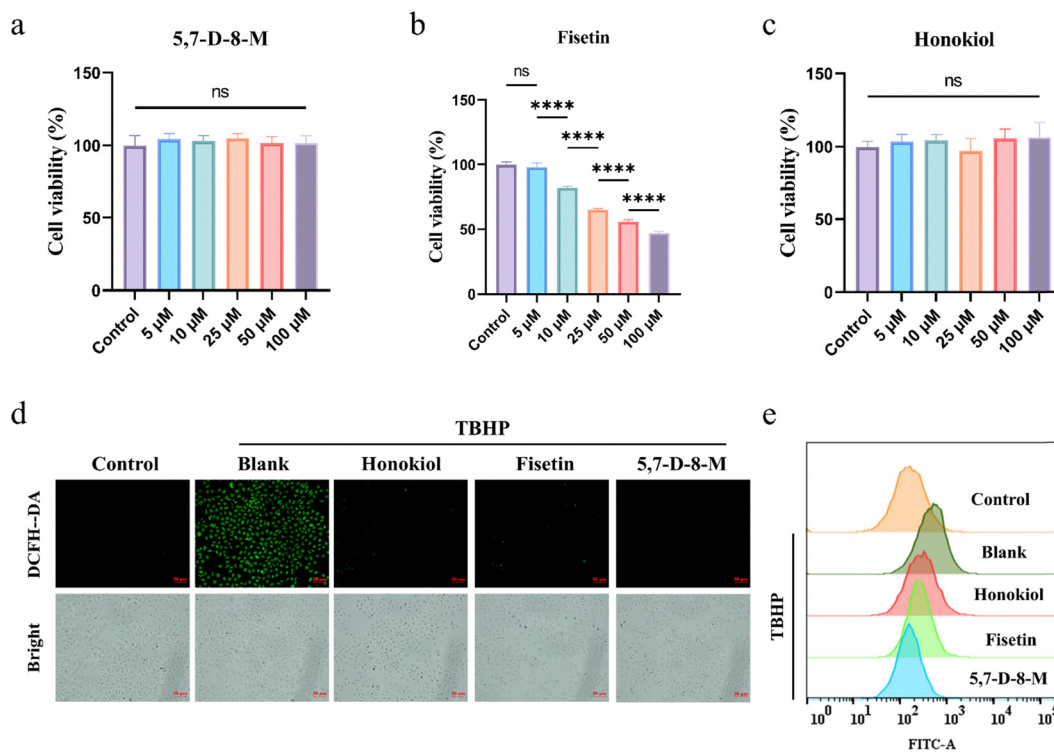




**Fig. 3** Molecular docking results of antioxidant-related proteins NRF2 and HO-1 with model screening compounds. Predicted binding mode of 5,7-diacetoxy-8-methoxyflavone, fisetin and honokiol docked with the (a)–(c) NRF2 and (d)–(f) HO-1, respectively. (Bonding residues as yellow bars, critical hydrogen bonds as yellow dashed lines).

stable molecular interactions within the catalytic pocket of heme oxygenase-1 (HO-1), as evidenced by computational docking analyses. Molecular docking results evaluating the interaction

between the screened compounds and superoxide dismutase (SOD) are provided in Fig. S2a–c (ESI<sup>†</sup>). The above results show the potential binding of 5,7-diacetoxy-8-methoxyflavone, fisetin



**Fig. 4** *In vitro* antioxidant assessment of model screening compounds. (a)–(c) Cell viability after HUVECs incubating with 5  $\mu$ M 5,7-diacetoxy-8-methoxyflavone, fisetin and honokiol for 24 h, respectively. TBHP-stimulated HUVECs after 24 h of treatment with candidate compounds (d) fluorescence images of ROS, (e) flow cytometry signal changes.



and honokiol to antioxidant target molecules. The antioxidant properties of the compounds screened for specific models still need to be systematically evaluated in biological experiments (Fig. 3).

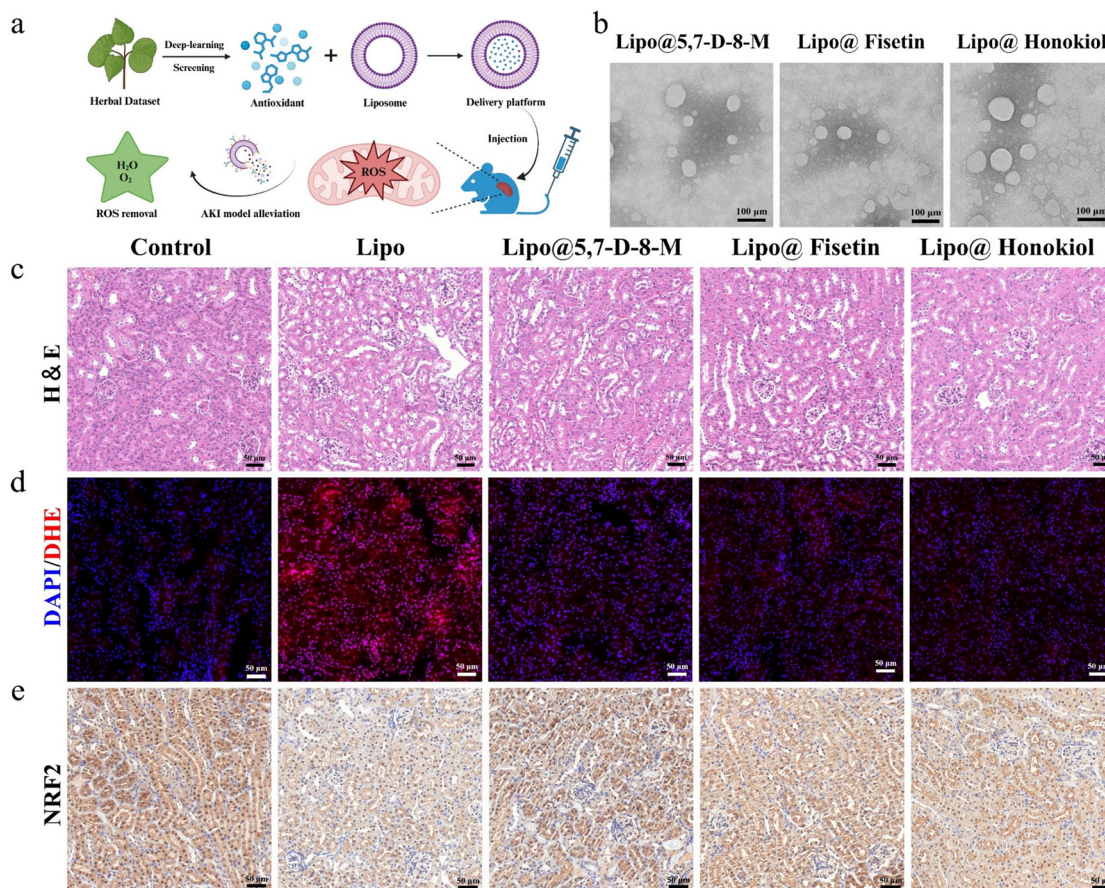
### 2.3 Verification of antioxidant properties *in vitro*

Considering candidate compounds potential value for *in vivo* antioxidant therapy, we first evaluated the cytotoxicity to HUVECs using the CCK-8 method and screened for suitable concentrations of action. The results showed that 5,7-diacetoxy-8-methoxyflavone and honokiol had no significant effect on cell proliferation ( $P > 0.05$ ) in the concentration range of 5–100  $\mu\text{M}$ , whereas 10  $\mu\text{M}$  fisetin significantly inhibited HUVECs proliferation ( $P < 0.001$ ) (Fig. 4a–c). In view of this result, and to ensure consistency in the treatment concentrations of the three compounds, 5  $\mu\text{M}$  was finally adopted for subsequent experiments. To verify the antioxidant effect, we used the ROS fluorescent probe combined with flow cytometry to detect the oxidative stress model of HUVEC induced by 0.2 mM TBHP. Compared with the positive control group, the intensity of green fluorescence in the cells was significantly reduced after 24 hours of treatment with 5  $\mu\text{M}$  5,7-diacetoxy-8-methoxyflavone, fisetin and

honokiol (Fig. 4d), and this trend was further confirmed by quantitative analysis using flow cytometry (Fig. 4e). The above results show that the candidate compounds have significant antioxidant activity *in vitro*.

### 2.4 Therapeutic effect of antioxidant liposome delivery platform *in vivo*

Based on the solubility differences of the three candidate compounds as well as their off-target and rapid scavenging effects *in vivo*, we further constructed functional liposomal drug platforms to improve their *in vivo* bioavailability and systematically evaluated their *in vivo* antioxidant effects (Fig. 5a). Transmission electron microscopy (TEM) results showed that the lipo@5,7-diacetoxy-8-methoxyflavone, lipo@fisetin, and lipo@honokiol group all showed typical phospholipid bilayer structure with an average diameter of about 80 nm (Fig. 5b). Tissue organ damage after acute ischemia-reperfusion is mainly caused by ROS excess, so we constructed a mouse kidney ischemia-reperfusion model to evaluate the *in vivo* antioxidant activity of antioxidant liposomes.<sup>21</sup> H&E staining results displayed that the liposome group showed histological features typical of renal tubular damage, with obvious dilated renal tubules, extensive cell debris and tubular casts, and



**Fig. 5** *In vivo* antioxidant assessment of model screening compounds. (a) Schematic diagram of the *in vivo* experimental process. (b) Observation of liposome-loaded compounds under a transmission electron microscope. Scale bar = 100  $\mu\text{m}$ . (c) Kidney H&E staining of mice with different treatments. Scale bar = 50  $\mu\text{m}$ . (d) DHE staining in mice with different treatments. Scale bar = 50  $\mu\text{m}$ . (e) Immunohistochemical staining with NRF2 in mice with different treatments. Scale bar = 50  $\mu\text{m}$ .



atrophied and deformed glomeruli. In contrast, after treatment with lipo@5,7-diacetoxy-8-methoxyflavone, lipo@fisetin and lipo@honokiol, renal tubular damage was effectively alleviated and there was less cellular debris and casts in the lumen. At the same time, the glomerulus remained relatively intact (Fig. 5c). Dihydroethidium (DHE) is a widely used redox-sensitive fluorescent probe that is specific for ROS such as superoxide and hydrogen peroxide. After treatment with the three antioxidant liposomes, the ROS level in the kidney was significantly lower than that in the empty liposome group (Fig. 5d). Nrf2 is a key regulator of oxidative stress response. Increased expression of Nrf2 suggests that drug-loaded liposomes may exert a nephroprotective effect by activating endogenous antioxidant defense pathways. Immunohistochemical analysis showed that the three drug-loaded liposome treatment groups significantly upregulated the expression of Nrf2 in the renal tissue of the acute kidney injury model compared to the empty liposome control group (Fig. 5e). The above results indicate that antioxidant liposomes derived from natural herbs still maintain significant antioxidant effects *in vivo* and have improved bioavailability.

### 3. Conclusion

In this study, we developed a two-step BERT-based framework model that efficiently screens antioxidant compounds from natural herbal data through pre-training and fine-tuning. Additionally, we constructed a versatile antioxidant liposome delivery platform based on the physicochemical properties of the candidate compounds, achieving effective antioxidant activity both *in vitro* and *in vivo*. This model demonstrates a performance and accuracy improvement of approximately 20% compared to the traditional machine learning approaches, significantly accelerating the discovery of antioxidant compounds from natural herbs. Through molecular docking and biological experiments, we successfully validated the excellent antioxidant effects of the compounds identified by our model. Furthermore, we integrated functional liposome modification technology to effectively encapsulate these candidate compounds, resulting in successful treatment in an animal model of ischemic acute kidney injury (AKI). Our study confirms the efficacy of deep learning in expediting the screening of antioxidant compounds and proposes a synergistic strategy that combines computational screening with a liposomal delivery platform, offering a novel paradigm for investigating oxidative stress-related diseases and the development of functional liposomes.

## 4. Materials and methods

### 4.1 Dataset curation

Our dataset for fine-tuning consists of 584 compounds with antioxidant activity and 983 compounds with oxidizing activity, as previously validated through experimental research. Each compound in the training set was assigned a binary label—1 for antioxidant activity and 0 for oxidizing activity, to fine-tune the pre-trained BERT model for a classification task aimed at screening potential antioxidants. Additionally, a dataset of

2882 natural herbs have been collected, which will be analyzed using the fine-tuned model to identify promising antioxidant candidates. All SMILES representations across training and test datasets were standardized to canonical form using RDKit.<sup>22</sup>

### 4.2 Model architecture and pre-training

In this work, we employed MTL-BERT-MEDIUM<sup>23</sup> as the foundation model, augmented with task-specific classification layers. The architecture comprises 8 Transformer layers with 8 attention heads and 256 hidden dimensions per layer, following standard BERT<sup>22</sup> pretraining procedures. BERT is a neural network architecture that learns contextual relationships by processing entire sequences of tokens simultaneously. Its core mechanism is multi-head self-attention, which captures long-range dependencies more effectively than traditional recurrent models.<sup>24</sup> In the self-attention operation, each token in the input is transformed into three vectors—query, key, and value—and combined as:

$$Z = \phi \left( \frac{(XW^Q)(XW^K)^T}{\sqrt{d_k}} \right) XW^V$$

where  $X$  is the input feature matrix,  $\phi$  is the softmax function,  $W^Q$ ,  $W^K$  and  $W^V$  are learnable weight matrices, and  $d_k$  is a scaling factor. The resulting matrix  $Z$  represents the attended output, which is then processed by feed-forward layers with residual connections and layer normalization. BERT employs multiple attention heads in each layer to learn different aspects of the input sequence.

Here, we employed the MTL-BERT-MEDIUM architecture and conducted our own pre-training on the ChEMBL\_v35 dataset. According to their pre-training experience settings, the learning rate was set to  $1 \times 10^{-4}$ , and the batch size was set to 512. The pre-training process was terminated after 40 epochs as additional training yielded marginal performance enhancements. Following the MTL-BERT-MEDIUM pre-training strategy, we only employed the masked language modeling (MLM) task, unlike the original BERT approach. SMILES strings lack the sequential narrative structure found in natural language, where sentence order and relationships are meaningful. Previous research<sup>25</sup> has demonstrated that effective language models can be developed without relying on inter-sequence relationships. Consequently, we focused solely on the masked token recovery task.

### 4.3 Model finetuning and predictions

During the fine-tuning stage, the pre-trained BERT model is used to screen potential antioxidant compounds. Each compound is represented by a tokenized SMILES string, which serves as the input to the model. For binary classification, the representation at the first position of the encoder output is extracted and mapped to a single logit score using a two-layer multi-layer perceptron (MLP):

$$\hat{y} = W_2 \cdot \text{LeakyReLU}(\text{Dropout}(W_1 Z^{[p1]} + b_1)) + b_2$$

Here,  $Z^{[p1]}$  represents the encoded representation for classification task, which first undergoes a linear projection to an expanded dimension, parameterized by  $W_1$  and bias  $b_1$ ,



followed by dropout regularization, then a LeakyReLU activation. The transformed representation is then passed through a second linear transformation back to a single dimension. The output  $\hat{y}$  represents the raw logit score.

To optimize the model, the encoder is fine-tuned using BCEWithLogitsLoss.

#### 4.4 Model interpretability

With the help of self-attention mechanism of BERT, we can quantify token importance in BERT by analyzing attention patterns from the classification token in the final transformer layer. The attention score can reveal which input tokens most significantly influence the model's representation when processing molecular data.

Our method extracts attention scores from BERT's final layer, where the model has developed its refined contextual understanding. Mathematically, if we denote the attention matrix of the last layer as  $A^L \in \mathbb{R}^{h \times n \times n}$ , the classification task-specific [p1] token attention scores are:

$$a_{p1} = \frac{1}{h} \sum_{i=1}^h A_{i,0,:}^L$$

Here,  $A^L$  represents the attention matrix from the last ( $L$ -th) transformer layer, with dimensions corresponding to the number of attention heads ( $h$ ) and sequence length ( $n$ ). Each element  $A_{i,j,k}^L$  quantifies how much the  $i$ -th attention head at position  $j$  attends to position  $k$  in the sequence.

For molecular applications, we render molecules from SMILES notation using RDKit and overlay normalized attention values as color gradients, with deeper red indicating higher model attention.

#### 4.5 Experiment setup

To ensure a rigorous evaluation, five-fold cross-validation is employed. Given the inherent class imbalance in our dataset, we evaluated multiple strategies to address this challenge, including no adjustment, focal loss, up-sampling, and down-sampling (detailed comparison in Table S1, ESI<sup>†</sup>). Based on empirical performance, we implemented up-sampling of the minority class within each training fold. In each fold, the dataset is split that 80% of the data is used for training with up-sampling applied to balance the classes, while the remaining 20% was reserved for validation. This process is repeated across all five folds, ensuring that each sample was used for validation exactly once. Model training is conducted for 80 epochs per fold using the Adam optimizer. The learning rate is set to  $1 \times 10^{-5}$ , and the batch size is fixed at 64 to optimize computational efficiency and predictive performance. Performance evaluation is conducted using the area under the receiver operating characteristic curve (ROC-AUC) and accuracy.

After training, the five independently trained models are applied to a natural herb dataset to predict antioxidant probability scores. The final probability for each compound is obtained by averaging the outputs from all five models. Based on these aggregated probabilities, compounds are ranked to

facilitate the identification of those with the highest predicted antioxidant potential. All experiments here is conducted in Python 3.8 and Pytorch.<sup>26</sup>

#### 4.6 Molecular docking

Molecular docking was performed to test the ability of the compounds screened by the model to bind to antioxidant molecules. Nrf2, SOD, and HO-1 were selected as receptors and the compounds screened by the model were selected as ligands. The molecular structures of the receptors were downloaded from the protein database (<https://www.rcsb.org>). Screening criteria: (1) protein source organism: Homo sapiens; (2) refinement resolution  $< 2.5$  Å; (3) complete protein structure with corresponding ligand. The molecular structure of the ligand was downloaded from the PubChem database (<https://pubchem.ncbi.nlm.nih.gov>). The receptor protein was standardized using PyMOL software to remove water molecules and impurities, and then the original ligand was separated to obtain a standardized receptor. The ligand and receptor were then imported into the AutoDockTools software. Polar hydrogen atoms and Gasteiger charges were added to the receptor. The grid box was then manually adjusted using the grid tool until the receptor was fully enveloped. Finally, molecular docking was performed using the AutoDock Vina software, and the docking results were visualized and hydrogen bond formation evaluated using the PyMOL software.

### Authors statement

X. Zhang: conceptualization, investigation, methodology, visualization, writing – original draft. Z. Zheng: conceptualization, methodology, visualization, writing – review & editing. L. Xie: investigation, methodology, visualization, data curation. M. Yang: formal analysis, investigation, methodology. J. Wang: investigation, methodology, visualization. W. Wang: formal analysis, validation. S. Han: conceptualization, supervision, investigation, methodology. Z. Zhang: conceptualization, supervision, writing – original draft, writing – review & editing. J. Wu: conceptualization, funding acquisition, project administration, supervision, writing – review & editing.

### Conflicts of interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### Data availability

The data supporting this article have been included as part of the ESI.<sup>†</sup>

### Acknowledgements

This work was supported by the National Natural Science Foundation of China (No. 52173150 and U22A20315), the Open Research



Funds from the Sixth Affiliated Hospital of Guangzhou Medical University, Qingyuan People's Hospital (202301-211), Guangzhou Science and Technology Program City-University Joint Funding (No. 2023A03J0001 and 2024A03J0604), the Natural Science Foundation and Science and Technology Project of Guangdong Province (2022A1515010071 and 2021A1515012321).

## References

- 1 S. Shadfar, S. Parakh, M. S. Jamali and J. D. Atkin, *Transl. Neurodegener.*, 2023, **12**, 18.
- 2 A. C. Boese and S. Kang, *Redox Biol.*, 2021, **42**, 101870.
- 3 A. Daiber, O. Hahad, I. Andreadou, S. Steven, S. Daub and T. Münzel, *Redox Biol.*, 2021, **42**, 101875.
- 4 M. Mirahmad, S. Mohseni, O. Tabatabaei-Malazy, F. Esmaeili, S. Alatab, R. Bahramsoltani, H.-S. Ejtahed, H. Qulami, Z. Bitarafan, B. Arjmand and E. Nazeri, *Phyto-medicine*, 2023, **109**, 154615.
- 5 H. Wang, Z. Xu, X. Li, J. Sun, D. Yao, H. Jiang, T. Zhou, Y. Liu, J. Li, C. Wang, W. Wang and R. Yue, *Carbohydr. Polym.*, 2017, **176**, 99–106.
- 6 Y. Zhang, L.-H. Liu, B. Xu, Z. Zhang, M. Yang, Y. He, J. Chen, Y. Zhang, Y. Hu, X. Chen, Z. Sun, Q. Ge, S. Wu, W. Lei, K. Li, H. Cui, G. Yang, X. Zhao, M. Wang, J. Xia, Z. Cao, A. Jiang and Y.-R. Wu, *Acta Pharm. Sin. B*, 2024, **14**, 3476–3492.
- 7 P. Zhang, X. Wang, X. Cen, Q. Zhang, Y. Fu, Y. Mei, X. Wang, R. Wang, J. Wang, H. Ouyang, T. Liang, H. Xia, X. Han and G. Guo, *Natl. Sci. Rev.*, 2024, **12**, nwae451.
- 8 T. Wu, R. Lin, P. Cui, J. Yong, H. Yu and Z. Li, *J. Pharm. Anal.*, 2024, **14**, 101022.
- 9 F. Wong, S. Omori, N. M. Donghia, E. J. Zheng and J. J. Collins, *Nat. Aging*, 2023, **3**, 734–750.
- 10 J. M. Stokes, K. Yang, K. Swanson, W. Jin, A. Cubillos-Ruiz, N. M. Donghia, C. R. MacNair, S. French, L. A. Carfrae, Z. Bloom-Ackermann, V. M. Tran, A. Chiappino-Pepe, A. H. Badran, I. W. Andrews, E. J. Chory, G. M. Church, E. D. Brown, T. S. Jaakkola, R. Barzilay and J. J. Collins, *Cell*, 2020, **180**, 688–702.e13.
- 11 E.-A. Kim, H. G. Choi, B. L. Nguyen, S.-J. Oh, S.-B. Lee, S. H. Bae, S. Y. Park, J. O. Kim, S. H. Kim and S.-J. Lim, *J. Controlled Release*, 2024, **366**, 410–424.
- 12 M. Dymek and E. Sikora, *Adv. Colloid Interface Sci.*, 2022, **309**, 102757.
- 13 J. Di, F. Xie and Y. Xu, *Adv. Drug Delivery Rev.*, 2020, **154–155**, 151–162.
- 14 D. E. Large, R. G. Abdelmessih, E. A. Fink and D. T. Auguste, *Adv. Drug Delivery Rev.*, 2021, **176**, 113851.
- 15 Z.-C. Sun, R. Liao, C. Xian, R. Lin, L. Wang, Y. Fang, Z. Zhang, Y. Liu and J. Wu, *J. Controlled Release*, 2024, **375**, 300–315.
- 16 M. Hirohara, Y. Saito, Y. Koda, K. Sato and Y. Sakakibara, *BMC Bioinf.*, 2018, **19**, 83–94.
- 17 L. Huang, S. Jia, R. Wu, Y. Chen, S. Ding, C. Dai and R. He, *Food Chem.*, 2022, **396**, 133713.
- 18 J. Zhang, Y. Li, J. Wan, M. Zhang, C. Li and J. Lin, *Phytomedicine*, 2022, **104**, 154259.
- 19 M. Fernanda Arias-Santé, J. Fuentes, C. Ojeda, M. Aranda, E. Pastene and H. Speisky, *Food Chem.*, 2024, **435**, 137487.
- 20 Y. Zhou, J. Tang, J. Lan, Y. Zhang, H. Wang, Q. Chen, Y. Kang, Y. Sun, X. Feng, L. Wu, H. Jin, S. Chen and Y. Peng, *Acta Pharm. Sin. B*, 2023, **13**, 577–597.
- 21 F. Zeng, Y. Qin, S. Nijati, Y. Liu, J. Ye, H. Shen, J. Cai, H. Xiong, C. Shi, L. Tang, C. Yu and Z. Zhou, *Adv. Sci.*, 2024, **11**, 2403305.
- 22 G. Landrum, *Release*, 2013, **1**, 4.
- 23 X.-C. Zhang, C.-K. Wu, J.-C. Yi, X.-X. Zeng, C.-Q. Yang, A.-P. Lu, T.-J. Hou and D.-S. Cao, *Research*, 2022, **2022**, 0004.
- 24 V. Ashish, *Advances in neural information processing systems*, 2017, **30**, 1.
- 25 Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer and V. Stoyanov, RoBERTa: A Robustly Optimized BERT Pretraining Approach, *arXiv*, 2019, preprint, arXiv:1907.11692, DOI: [10.48550/arXiv.1907.11692](https://doi.org/10.48550/arXiv.1907.11692).
- 26 A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein and L. Antiga, PyTorch: An Imperative Style, High-Performance Deep Learning Library, *arXiv*, 2019, preprint, arXiv:1912.01703, DOI: [10.48550/arXiv.1912.01703](https://doi.org/10.48550/arXiv.1912.01703).

